

## Joint Entropy and Conditional Entropy

Consider a vector valued variable  $(X, Y)$ , The joint entropy  $H(X, Y)$  is

$$H(X, Y) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x, y)$$

$$H(X, Y) = -E\{\log p(X, Y)\}$$

The conditional entropy is defined as

$$\begin{aligned} H(Y|X) &= \sum_{x \in X} p(x) H(Y|X=x) \\ &= - \sum_{x \in X} p(x) \sum_{y \in Y} p(y|x) \log p(y|x) \\ &= - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(y|x) \end{aligned}$$

Chain rule:

$$H(X, Y) = H(X) + H(Y|X)$$

$$\begin{aligned} H(X, Y) &= - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x, y) \\ &= - \sum_{x \in X} \sum_{y \in Y} p(x, y) (\log p(x) + \log p(y|x)) \\ &= - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x) - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(y|x) \\ &= - \sum_{x \in X} p(x) \log p(x) - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(y|x) \\ &= H(X) + H(Y|X) \end{aligned}$$

Corollary:  $H(X, Y|Z) = H(X|Z) + H(Y|X, Z)$  - show!

Example:

X \ Y	1	2	3	4	P(Y)
1	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{32}$	$\frac{1}{32}$	$\frac{1}{4}$
2	$\frac{1}{16}$	$\frac{1}{8}$	$\frac{1}{32}$	$\frac{1}{32}$	$\frac{1}{4}$
3	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{4}$
4	$\frac{1}{4}$	0	0	0	$\frac{1}{4}$
P(X)	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$	

$$H(X) = \frac{7}{4} \text{ SL bits}$$

$$H(Y) = 2 \text{ SL bits}$$

$$H(X|Y) = \sum_{i=1}^4 P(Y=i) H(X|Y=i)$$

$$\begin{aligned} H(X|Y) &= \sum_{x \in X} \sum_{y \in Y} P(x,y) (\log P(x|y)) \\ &= \frac{11}{8} \text{ bits} \end{aligned}$$

$$\begin{aligned} P(X=1|Y=1) \cdot P(Y=1) &= P(X=1, Y=1) \\ P(X=1|Y=1) &= \frac{P(X=1, Y=1)}{P(Y=1)} = \frac{\frac{1}{8}}{\frac{1}{4}} = \frac{1}{2} \end{aligned}$$

## Mutual Information

Consider two random variables  $X$  and  $Y$  with a joint probability distribution  $P(x, y)$  and marginal distributions  $p(x)$  and  $p(y)$ .

The mutual information  $I(X; Y)$  is relative entropy between the joint distribution and the product distribution  $p(x)p(y)$ , i.e.,

$$I(X; Y) = \sum_{x \in X} \sum_{y \in Y} P(x, y) \cdot \log \frac{P(x, y)}{P(x)P(y)}$$

We can show that  $I(X; Y) = H(X) - H(X|Y)$

$$\begin{aligned} I(X; Y) &= \sum_{x \in X} \sum_{y \in Y} P(x, y) \log \frac{P(x, y)}{P(x)P(y)} \\ &= \sum_{x \in X} \sum_{y \in Y} P(x, y) \log \frac{P(x|y)}{P(x)} \\ &= - \sum_{x \in X} \sum_{y \in Y} P(x, y) \log P(x) + \sum_{x \in X} \sum_{y \in Y} P(x, y) \log P(x|y) \\ &= - \sum_{x \in X} P(x) \log P(x) - \left( - \sum_{x \in X} \sum_{y \in Y} P(x, y) \log P(x|y) \right) \\ &= H(X) - H(X|Y) \end{aligned}$$

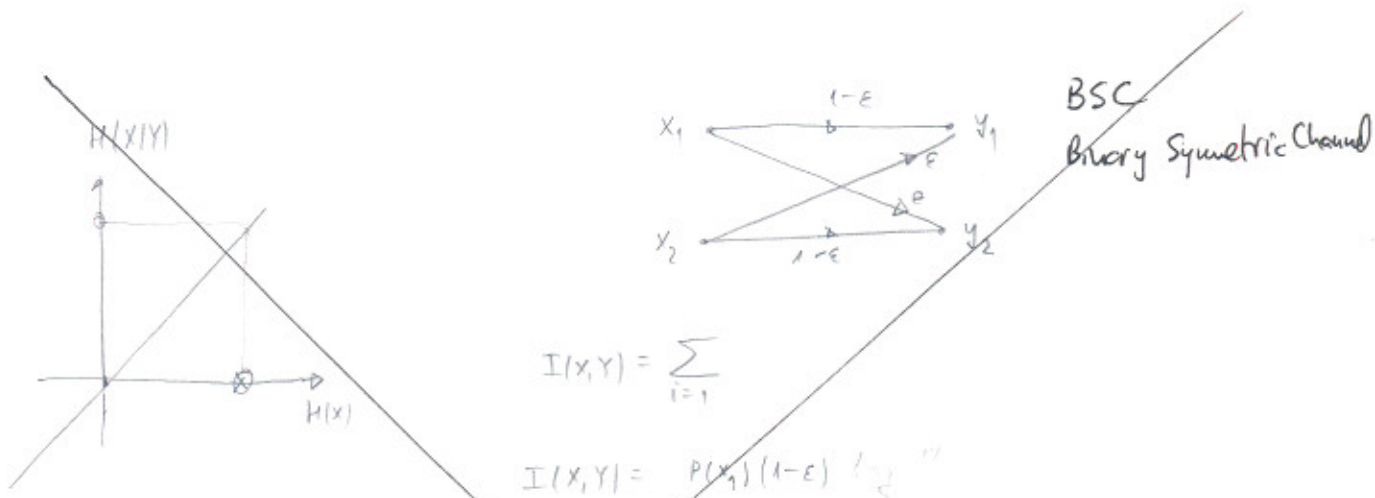
— The mutual information  $I(X; Y)$  is the reduction in the uncertainty of  $X$  due to the knowledge of  $Y$ .

— It follows

$$I(X; Y) = H(Y) - H(Y|X)$$

$I(X; X) = H(X) - H(X)$  is referred to as self information

$$\begin{aligned}
 I(X; Y) &= \sum_{i=1}^m \sum_{j=1}^n P(x_i) P(y_j | x_i) \log \frac{P(y_j | x_i)}{P(y_j)} \\
 &= \sum_{i=1}^m \sum_{j=1}^n P(x_i) P(y_j | x_i) (\log P(y_j | x_i) + \log P(y_j)) + \sum_{i=1}^m \sum_{j=1}^n P(x_i) P(y_j | x_i) (\log P(y_j)) \\
 &= \sum_{i=1}^m P(x_i) \sum_{j=1}^n P(y_j | x_i) \log P(y_j | x_i) + \sum_{j=1}^n (\log P(y_j)) \underbrace{\sum_{i=1}^m P(x_i) P(y_j | x_i)}_{P(y_j)} \\
 &= -H(Y|X) + H(Y)
 \end{aligned}$$



$$I(X, Y) = \sum_{i=1}^m$$

$$I(X, Y) = P(x_1) (1-\epsilon) \log \dots$$

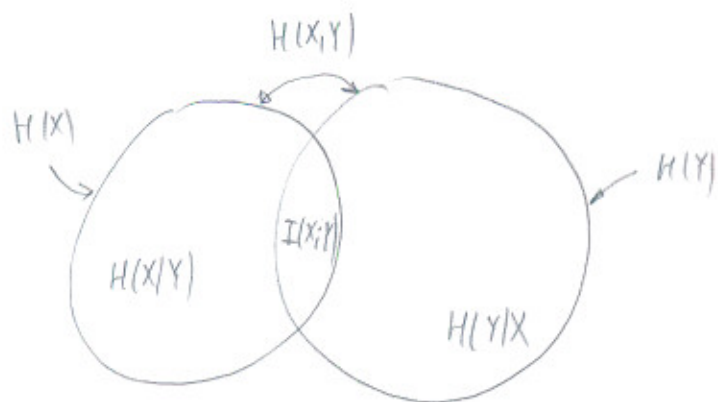
$$I(X, Y) = H(Y) - H(Y|X)$$

$$H(Y|X) = \sum_{i=1}^m \sum_{j=1}^n P(x_i, y_j) \log \frac{1}{P(y_j | x_i)}$$

$$(1-p)(1-\epsilon) \log \frac{1}{1-\epsilon} + (1-p)\epsilon \log \frac{1}{\epsilon}$$

$$p(1-\epsilon) \log \frac{1}{1-\epsilon} + p\epsilon \log \frac{1}{\epsilon}$$

$$(1-\epsilon) \log \frac{1}{1-\epsilon} + \epsilon \log \frac{1}{\epsilon}$$



$$I(X;Y) = H(X) - H(X|Y)$$

$$I(X;Y) = H(X) + H(Y) - H(X,Y)$$

Some interesting properties:

$$H(X_1, X_2, \dots, X_n) = H(X_1) + H(X_2|X_1) + H(X_3|X_1, X_2) + \dots + H(X_n|X_1, X_2, \dots, X_{n-1})$$

$$I(X_1, X_2, \dots, X_n; Y) = I(X_1; Y) + I(X_2; Y|X_1) + I(X_3; Y|X_1, X_2) + \dots + I(X_n; Y|X_1, X_2, \dots, X_{n-1})$$

$$I(X; Y) \geq 0$$

$$H(X|Y) \leq H(X)$$

$$H(X_1, X_2, \dots, X_n) \leq \sum_{i=1}^n H(X_i)$$

Data processing inequality

If  $X, Y, Z$  form a Markov chain ( $P(X, Y, Z) = P(X)P(Y|X)P(Z|Y)$ ) then

$$I(X; Y) \geq I(X; Z)$$

$$\stackrel{\geq 0}{I(X; Y, Z)} = I(X; Z) + \overbrace{I(X; Y|Z)}^{\geq 0}$$

$$\text{or } = I(X; Y) + I(X; Z|Y)$$

$$I(X; Y) \geq I(X; Z) \quad \leftarrow$$

- since  $X$  and  $Z$  are conditionally independent given  $Y$   $I(X; Z|Y) = 0$
- since  $I(X; Y|Z) \geq 0$

$$I_{X;Y}(x_i; y_j) = \log \frac{P_{XY}(x_i, y_j)}{P_X(x_i)}$$

$$I_{X;Y}(x_i; y_j) = \log \frac{P_{XY}(x_i, y_j) P(Y_j)}{P_X(x_i) P(Y_j)}$$

$$= \log \frac{P_{XY}(x_i, y_j)}{P_X(x_i) P_Y(y_j)} = \log \frac{P_{Y|X}(y_j | x_i) P(x_i)}{P_X(x_i) P_Y(y_j)}$$

$$= \log \frac{P_{Y|X}(y_j | x_i)}{P_Y(y_j)}$$

$$= I_{Y;X}(y_j; x_i)$$

$$I(X; Y) = \sum_{i=1}^I \sum_{j=1}^J P_{XY}(x_i, y_j) \log \frac{P_{XY}(x_i, y_j)}{P_X(x_i) P_Y(y_j)}$$

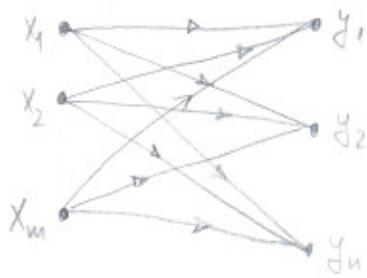
$$I(u_1; u_2 | u_3) = \log \frac{P(u_1 | u_2, u_3)}{P(u_1 | u_3)} = I(u_1 | u_3) - I(u_1 | u_2, u_3)$$

$$I(u_1; u_2 | u_3) = \sum_{u_1} \sum_{u_2} \sum_{u_3} P(u_1, u_2, u_3) \cdot \log \frac{P(u_1 | u_2, u_3)}{P(u_1 | u_3)}$$

$$= H(u_1 | u_3) - H(u_1 | u_2, u_3)$$



# Average Mutual Information



Discrete  
Memoryless  
Channel

$$P(x_i) \quad 1 \leq i \leq m$$

$$P(y_j) \quad 1 \leq j \leq n$$

$$P(y_j | x_i)$$

$$H(X) = \sum_{i=1}^m P(x_i) \log \frac{1}{P(x_i)}$$

$$H(Y) = \sum_{j=1}^n P(y_j) \log \frac{1}{P(y_j)}$$

$$I(x_i; y_j) = \log \frac{P(x_i, y_j)}{P(x_i)P(y_j)}$$

← mutual information  
- information provided about  $x_i$  by  $y_j$

$$I(x_i; y_j) = \log \frac{P(x_i, y_j)}{P(x_i)P(y_j)}$$

$$H(X|Y) = \sum_{i=1}^m \sum_{j=1}^n P(x_i, y_j) \log \frac{1}{P(x_i|y_j)}$$

← amount of average uncertainty remaining in  $X$  after observation of  $Y$

$$I(X; Y) = \sum_{i=1}^m \sum_{j=1}^n P(x_i, y_j) I(x_i; y_j)$$

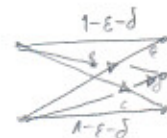
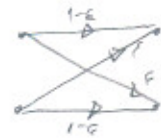
← average mutual information

$$I(X; Y) = H(X) - H(X|Y)$$

$$I(X; Y) = H(Y) - H(Y|X)$$

## Channel Models

Binary Symmetric Channel



Binary Erasure Channel

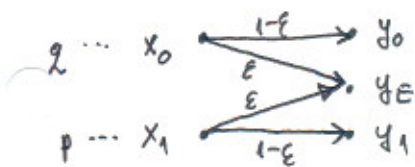
$$I(u_1; u_2 | u_3) = \sum_{u_1} \sum_{u_2} \sum_{u_3} P(u_1, u_2, u_3) \log \frac{P(u_1, u_2, u_3)}{P(u_1, u_2, u_3)}$$

$$I(X, Y) = \sum_{i=1}^m \sum_{j=1}^n P(x_i) P(y_j | x_i) \log \frac{P(y_j | x_i)}{P(y_j)}$$

$$= \sum_{i=1}^m \sum_{j=1}^n P(x_i, y_j) \log P(y_j | x_i) - \sum_{j=1}^n \log P(y_j) \sum_{i=1}^m P(x_i) P(y_j | x_i)$$

$$= -H(Y|X) + H(Y)$$

# BEC



$$P(x_0, y_0) = P(x_0) \cdot P(y_0 | x_0) = q(1-\epsilon)$$

$$P(x_0, y_E) = P(x_0) P(y_E | x_0) = q \cdot \epsilon$$

$$P(x_0, y_1) = P(x_0) \underbrace{P(y_1 | x_0)}_0 = 0$$

$$P(x_1, y_0) = 0$$

$$P(x_1, y_E) = p \epsilon$$

$$P(x_1, y_1) = p(1-\epsilon)$$

$$P(y_0) = \sum_i P(x_i) P(y_0 | x_i) = q(1-\epsilon)$$

$$P(y_E) = \sum_i P(x_i) P(y_E | x_i) = q\epsilon + p\epsilon = \epsilon$$

$$P(y_1) = \sum_i P(x_i) P(y_1 | x_i) = p(1-\epsilon)$$

$$I(X; Y) = H(Y) - H(Y|X)$$

$$H(Y) = -\sum_j P(y_j) \log P(y_j) = -q(1-\epsilon) \log q(1-\epsilon) - \epsilon \log \epsilon - p(1-\epsilon) \log p(1-\epsilon)$$

$$H(Y, X) = -\sum_i \sum_j P(x_i, y_j) \log P(y_j | x_i)$$

$$= -q(1-\epsilon) \cdot \log(1-\epsilon) - q \cdot \epsilon \log \epsilon - p \cdot \epsilon \log \epsilon - p(1-\epsilon) \log(1-\epsilon)$$

$$\begin{aligned} H(Y) - H(Y, X) &= -q(1-\epsilon) \log q - q(1-\epsilon) \cancel{\log(1-\epsilon)} - p(1-\epsilon) \log p - p(1-\epsilon) \cancel{\log(1-\epsilon)} \\ &\quad + q(1-\epsilon) \cancel{\log(1-\epsilon)} + p(1-\epsilon) \cancel{\log(1-\epsilon)} \quad \begin{matrix} -\epsilon \log \epsilon \\ +\epsilon \log \epsilon \end{matrix} \\ &= -q(1-\epsilon) \log q - p(1-\epsilon) \log p \end{aligned}$$

$$I(X; Y) = -p(1-\epsilon) (q \log q - p \log p)$$

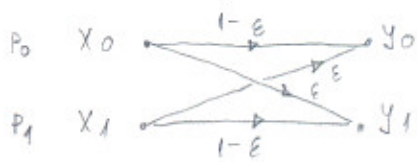
$$= (1-\epsilon) H(X) = H(X) - \epsilon H(X)$$

$$(H(X|Y) = \epsilon H(X))$$

$$I_{\max}(X; Y) = 1 - \epsilon$$



# Binary Symmetric Channel



$$I(X; Y) = H(Y) - H(Y|X)$$

$$P(Y_0) = P(X_0) \cdot (1 - \epsilon) + P(X_1) \cdot \epsilon = P(X_0) + (P(X_1) - P(X_0)) \epsilon$$

$$P(Y_1) = P(X_1) + (P(X_0) - P(X_1)) \epsilon$$

$$H(Y|X) = \sum_{i=0}^1 \sum_{j=0}^1 P(X_i, Y_j) \log \frac{1}{P(Y_j|X_i)} = \sum_{i=0}^1 \sum_{j=0}^1 P(X_i) P(Y_j|X_i) \log \frac{1}{P(Y_j|X_i)}$$

$$H(Y|X) = P_0(1 - \epsilon) \log \frac{1}{1 - \epsilon} + P_0 \cdot \epsilon \log \frac{1}{\epsilon} + P_1 \cdot \epsilon \log \frac{1}{\epsilon} + P_1(1 - \epsilon) \log \frac{1}{1 - \epsilon}$$

$$P_0 = P_1 = \frac{1}{2}$$

maximize

$$H(Y|X) = (1 - \epsilon) \log \frac{1}{1 - \epsilon} + \epsilon \log \frac{1}{\epsilon} = H(\epsilon)$$

$$H(Y) = 1$$

$$I(X; Y) = 1 - H(\epsilon)$$



How to choose input probabilities to maximize  $I(X, Y)$   
(DMC)

$$I(X_i; Y) = \sum_{j=1}^n P(y_j | X_i) \log \frac{P(y_j | X_i)}{P(y_j)}$$

$$I(X; Y) = \sum_{i=1}^m P(x_i) I(X_i; Y)$$

$$I(X_i; Y) = \begin{cases} = C & P(x_i) > 0 \\ \leq C & P(x_i) = 0 \end{cases}$$

Arimoto-Blahut algorithm

## Example: Symmetric channels

- Each row of  $\Pi$  contains the same set of probabilities  $\{p_j\}_{j=1}^n$  and each column contains the same set of numbers  $\{q_i\}_{i=1}^m$

- Examples:

$$\Pi = \begin{bmatrix} \frac{1}{2} & \frac{1}{3} & \frac{1}{6} \\ \frac{1}{6} & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{2} & \frac{1}{6} & \frac{1}{3} \end{bmatrix}$$

$$\Pi = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{6} & \frac{1}{6} \\ \frac{1}{6} & \frac{1}{6} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

- An important property of the symmetric channel is that  $H(Y|X)$  is independent of  $P(X_i)$  and depends only on the channel transition probability matrix.

$$H(Y|X) = - \sum_{i=1}^m \sum_{j=1}^n P(X_i) P(Y_j|X_i) \log \frac{1}{P(Y_j|X_i)}$$

$$= - \sum_{i=1}^m P(X_i) \underbrace{\sum_{j=1}^n P(Y_j|X_i) \log \frac{1}{P(Y_j|X_i)}}_{\text{independent on } i}$$

$$= - \underbrace{\sum_{i=1}^m P(X_i)}_{=1} \underbrace{\sum_{j=1}^n P_j \log \frac{1}{P_j}}_{\text{still independent on } i}$$

$$= - \sum_{j=1}^n P_j \log \frac{1}{P_j}$$

- We have proved that the conditional entropy  $H(Y|X)$  does not depend on the input probability distribution. Thus the problem of maximizing  $I(X;Y) = H(Y) - H(Y|X)$  reduces to the problem of maximizing the output entropy  $H(Y)$ .

- We know that  $H(Y) \leq \log_2 n$  where the equality is achieved when  $P(Y_j) = \frac{1}{n} \quad 1 \leq j \leq n$

- We prove that the output symbols are equally likely when so are input symbols

$$P(Y_j) = \sum_{i=1}^w P(X_i) P_{ij} = \sum_{i=1}^w P(X_i) \cdot 2_i \Rightarrow \text{if } P(X_i) = \frac{1}{w}$$

$\uparrow$   
 $\oplus$  independent on  $j$  -

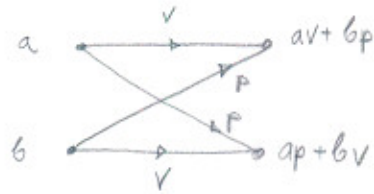
$$P(Y_j) = \frac{1}{w} \sum_{i=1}^w 2_i - \text{all equal}$$

- Thus all symbols  $y \in Y$  have the same probability and the capacity of a symmetric channel is given by

$$C = \log n + \sum_{j=1}^n P_j \log_2 \frac{1}{P_j}$$

Example

Binary Symmetric Channel (BSC)



$$I(X;Y) = H(Y) - H(Y|X)$$

$$= H(Y) - \left( a \cdot v \cdot \log \frac{1}{v} + a \cdot p \log \frac{1}{p} + b \cdot p \log \frac{1}{p} + b \cdot v \log \frac{1}{v} \right)$$

$$= H(Y) - \underbrace{(a+b) \left( p \log \frac{1}{p} + v \log \frac{1}{v} \right)}_{H(P)}$$

$$= H(av+bv) - H(P)$$

$$H(X) = x \log \frac{1}{x} + (1-x) \log \frac{1}{1-x}$$

$$\begin{cases} p \leq av+bv \leq v & \text{for } p \leq \frac{1}{2} \\ p \geq av+bv \geq v & \text{for } p \geq \frac{1}{2} \end{cases} \Rightarrow H(av+bv) \geq H(p)$$

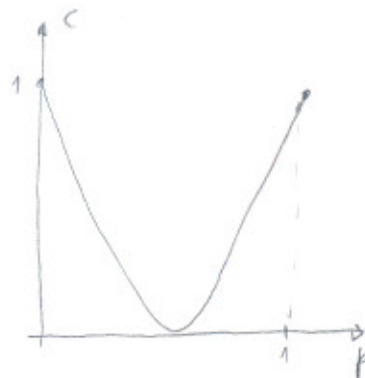
$$\Rightarrow I(X;Y) \geq 0$$

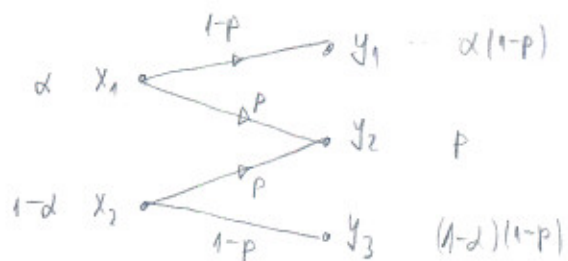
For a fixed  $p$ ,  $H(av+bv)$  achieves maximum

$$\text{for } a = \frac{1}{2} \quad H(av+bv) = H\left(\frac{1}{2}p + \frac{1}{2}(1-p)\right) = H\left(\frac{1}{2}\right) = 1 \quad \frac{\text{bit}}{\text{symbol}}$$

$$I_{\max}(X;Y) = 1 - H(p)$$

$$C = 1 - H(p)$$





$$I(X; Y) = H(Y) - H(Y|X)$$

$$= -\alpha(1-p) \log \alpha(1-p) - p \log p - (1-\alpha)(1-p) \log (1-\alpha)(1-p)$$

$$+ \alpha(1-p) \log (1-p) + \alpha p \log p + (1-\alpha)p \log p + (1-\alpha)(1-p) \log (1-p)$$

$\Rightarrow$

$$C = \max_{P(X;)} I(X; Y) = I(X; Y) \Big|_{\alpha = \frac{1}{2}} = 1-p$$

